

# Analysis of Resparsification

(draft)

Jakub Pachocki  
pachocki@cs.cmu.edu

## Abstract

We show that schemes for sparsifying matrices based on iteratively resampling rows yield guarantees matching classic ‘offline’ sparsifiers (see e.g. [SS08]). In particular, this gives a formal analysis of a scheme very similar to the one proposed by Kelner and Levin [KL13].

## 1 Introduction

In this note, we deal with the problem of spectrally approximating, or *sparsifying*, a tall  $n \times d$  matrix  $\mathbf{A}$ . Our goal will be to find an  $\mathcal{O}(d \log d \epsilon^{-2}) \times d$  matrix  $\tilde{\mathbf{A}}$  such that  $\|\tilde{\mathbf{A}}\mathbf{x}\|_2 \approx \|\mathbf{A}\mathbf{x}\|_2$  for all  $\mathbf{x}$ ; we will focus on the streaming setting, where we read rows of  $\mathbf{A}$  one by one. We will also require that our  $\tilde{\mathbf{A}}$  preserves the structure of  $\mathbf{A}$ : that is, it will only consist of a subset of (reweighted) rows of  $\mathbf{A}$ .

It is well known that sampling  $\mathcal{O}(d \log d \epsilon^{-2})$  rows of  $\mathbf{A}$  with probabilities proportional to their *leverage scores* yields a  $(1 + \epsilon)$ -spectral approximation to  $\mathbf{A}$  (see e.g. [SS08, LMP13, CLM<sup>+</sup>15, CMP16]).

Kelner and Levin [KL13] proposed a scheme for computing  $\tilde{\mathbf{A}}$  through simply storing the rows of  $\mathbf{A}$  in memory as they are read, and applying the above sparsification procedure whenever the matrix becomes too large. The main contribution of this note is a rigorous proof of such a scheme.

We believe the fact that the *resparsification* paradigm yields algorithms with guarantees as good as offline sparsifiers has implications reaching beyond the streaming setting; similar schemes have been used for example in designing fast Laplacian solvers [PS14, KS16].

## 2 Preliminaries

Let  $\mathbf{a}_1, \dots, \mathbf{a}_n$  be vectors in  $\mathbb{R}^d$ , such that

$$\sum_{i=1}^n \mathbf{a}_i \mathbf{a}_i^\top = \mathbf{I}.$$

We will use the following two theorems due to Tropp [Tro11, Tro12]:

**Theorem 2.1 (Matrix Chernoff)** *Consider a finite sequence  $\{\mathbf{X}_k\}$  of independent, random, self-adjoint matrices with dimension  $d$ . Assume that each random matrix satisfies*

$$\mathbf{X}_k \succeq \mathbf{0} \text{ and } \|\mathbf{X}_k\|_2 \leq R \text{ almost surely.}$$

---

Define

$$\mu_{\max} := \left\| \sum_k \mathbb{E} [\mathbf{X}_k] \right\|_2.$$

Then

$$\mathbf{P} \left[ \left\| \sum_k \mathbf{X}_k \right\|_2 \geq (1 + \delta) \mu_{\max} \right] \leq d \cdot \left( \frac{e^\delta}{(1 + \delta)^{1+\delta}} \right)^{\mu_{\max}/R} \text{ for } \delta \geq 0.$$

**Theorem 2.2 (Matrix Freedman)** *Let  $\mathbf{Y}_0, \mathbf{Y}_1, \mathbf{Y}_2, \dots$  be a matrix martingale whose values are self-adjoint matrices with dimension  $d$ , and let  $\mathbf{X}_1, \mathbf{X}_2, \dots$  be the difference sequence. Assume that the difference sequence is uniformly bounded in the sense that*

$$\|\mathbf{X}_k\|_2 \leq R \text{ almost surely, for all } k.$$

*Define the predictable quadratic variation process of the martingale:*

$$\mathbf{W}_k := \sum_{j=1}^k \mathbb{E}_{j-1} [\mathbf{X}_j^2], \text{ for all } k.$$

*Then, for all  $t > 0$  and  $\sigma^2 > 0$ ,*

$$\mathbf{P} [\exists k : \|\mathbf{Y}_k\|_2 \geq t \text{ and } \|\mathbf{W}_k\|_2 \leq \sigma^2] \leq d \cdot \exp \left( -\frac{-t^2/2}{\sigma^2 + Rt/3} \right).$$

### 3 The Game

Fix  $\epsilon \in (0, \frac{1}{2})$ . Set  $c = 100\epsilon^{-2} \log d$ . Assume that for all  $i$ ,  $\mathbf{a}_i^\top \mathbf{a}_i \leq \frac{1}{c}$  (this assumption is technical and not actually necessary; larger rows are simply ignored in the game anyway). For  $i = 1, \dots, n$ , set  $w_i$  to 1.

We will analyze a game played by an adversary on the weights  $w_i$ . It consists of a single move, repeated while the game is not over:

1. The adversary picks any  $i$  such that  $w_i \neq 0$  and  $2w_i \mathbf{a}_i^\top \mathbf{a}_i \leq \frac{1}{c}$ .
2. Flip an unbiased coin. If it comes out heads, set  $w_i \leftarrow 2w_i$ ; otherwise, set  $w_i \leftarrow 0$ .

The game can end in one of two ways:

- The matrix  $\sum w_i \mathbf{a}_i \mathbf{a}_i^\top$  is not a  $(1 + \epsilon)$ -approximation to the identity; then, the adversary wins.
- The adversary has no more legal moves; then, the adversary loses.

We will show that with high probability, the adversary will not win the game.

**Theorem 3.1** *With high probability, the game defined above ends in a loss for the adversary.*

## 4 Bounding the Total Variation

For all  $i$ , let  $w'_i$  be the maximum value attained by  $w_i$  throughout the game.

**Lemma 4.1** *Whp, we have that*

$$\left\| \sum w_i'^2 (\mathbf{a}_i \mathbf{a}_i^\top)^2 \right\|_2 \leq \frac{4}{c}.$$

**Proof:** First of all, we consider a slightly modified version of the game: if the game ends in a victory for the adversary, we still let them perform legal moves while possible. This can only increase the  $w'_i$ s.

Now note that in this augmented game, the adversary's choices do not have any effect on the final values of  $w_i$  and  $w'_i$ s; for every  $i$ , the adversary will keep increasing  $w_i$  until it is either too large or 0. The  $w'_i$  are therefore independent.

Let  $\mathbf{X}'_i := w_i'^2 (\mathbf{a}_i \mathbf{a}_i^\top)^2$ , for  $i = 1, \dots, n$ . Our goal is now to bound the norm of the sum of the independent matrices  $\mathbf{X}'_i$ . First of all, note that we always have  $\|\mathbf{X}'_i\|_2 \leq \frac{1}{c^2}$ .

Let  $B_i$  be the maximum integer such that  $2^{B_i} \mathbf{a}_i \mathbf{a}_i^\top \leq \frac{1}{c}$ . We have

$$\begin{aligned} \mathbb{E} [\mathbf{X}'_i] &\preceq \sum_{k=0}^{B_i} 2^{-k} \cdot 4^k \cdot (\mathbf{a}_i \mathbf{a}_i^\top)^2 \\ &\preceq 2^{B_i+1} \cdot (\mathbf{a}_i \mathbf{a}_i^\top)^2 \\ &\preceq \frac{2}{c} \cdot \mathbf{a}_i \mathbf{a}_i^\top, \end{aligned}$$

and so

$$\left\| \sum_{i=1}^n \mathbb{E} [\mathbf{X}'_i] \right\|_2 \leq \frac{2}{c}.$$

The thesis follows from Theorem 2.1 with  $\delta = \frac{2}{c \cdot \mu_{\max}}$ ,  $R = \frac{1}{c^2}$ . ■

## 5 The Proof

We consider a martingale with the difference  $\mathbf{X}_j$  corresponding to the  $j$ -th move by the adversary, or  $\mathbf{0}$  if the adversary made less than  $j$  moves. Assume the adversary chooses row  $\mathbf{a}_i$ ; then we have

$$\mathbf{X}_j := \begin{cases} w_i \mathbf{a}_i \mathbf{a}_i^\top & \text{with probability } \frac{1}{2} \\ -w_i \mathbf{a}_i \mathbf{a}_i^\top & \text{otherwise.} \end{cases}$$

Let  $\{\mathbf{W}_k\}$  be the predictable quadratic variation process of the martingale given by the  $\{\mathbf{X}_j\}$ :

$$\mathbf{W}_k := \sum_{j=1}^k \mathbb{E}_{j-1} [\mathbf{X}_j^2].$$

---

**Lemma 5.1** *Whp, for all  $k$  we have that*

$$\|\mathbf{W}_k\|_2 \leq \frac{8}{c}.$$

**Proof:** We have

$$\mathbb{E}_{j-1} [\mathbf{X}_j^2] = w_i^2 (\mathbf{a}_i \mathbf{a}_i^\top)^2.$$

Therefore

$$\mathbf{W}_k \preceq \sum_{i=1}^n 2 \cdot w_i^2 (\mathbf{a}_i \mathbf{a}_i^\top)^2.$$

The thesis follows from Lemma 4.1. ■

**Proof of Theorem 3.1:** Note that  $\|\mathbf{X}_k\|_2 \leq \frac{1}{c}$ . Together with Lemma 5.1, applying Theorem 2.2 with  $t = \epsilon, \sigma^2 = \frac{8}{c}, R = \frac{1}{c}$  yields the thesis. ■

## 6 Application to Streaming Sparsification

Consider any sparsification algorithm that reads edges of the graph one by one and adds them to the sparsifier. At any time, the algorithm is free to take any edges whose leverage score is not too large, and remove them from the sparsifier with probability  $\frac{1}{2}$  and double their weight otherwise. Such an algorithm can be implemented in  $\mathcal{O}(d \log d \epsilon^{-2})$  space and nearly linear time (see Figure 1).

Theorem 3.1 gives that any such algorithm will output a sparsifier of the original graph whp. at the end. Some additional care is required to show that we can maintain a sparsifier to the current graph at all times.

$\tilde{\mathbf{A}} = \text{STREAMING-SAMPLE}(\mathbf{A}, \epsilon)$ , where  $\mathbf{A}$  is an  $n \times d$  matrix with rows  $\mathbf{a}_1, \dots, \mathbf{a}_n$ ,  $\epsilon \in (0, \frac{1}{2})$ .

1. Set  $c := 100 \log d / \epsilon^2$ .
2. Let  $\tilde{\mathbf{A}}$  be a  $0 \times d$  matrix.
3. For  $i = 1, \dots, n$ :
  - (a) Append  $\mathbf{a}_i$  to  $\tilde{\mathbf{A}}$ .
  - (b) If  $\tilde{\mathbf{A}}$  has more than  $20dc$  rows:
    - i. Let  $l_j$  be the leverage score of the  $j$ -th row of  $\tilde{\mathbf{A}}$ .
    - ii. While  $\tilde{\mathbf{A}}$  has more than  $10dc$  rows:
      - A. Pick an arbitrary row  $\tilde{\mathbf{a}}_j$  of  $\tilde{\mathbf{A}}$  with  $l_j$  less than  $\frac{1}{4c}$ .
      - B. With probability  $\frac{1}{2}$ , remove  $\tilde{\mathbf{a}}_j$  from  $\tilde{\mathbf{A}}$ ; otherwise, double  $\tilde{\mathbf{a}}_j$  and  $l_j$ .
4. Return  $\tilde{\mathbf{A}}$ .

Figure 1: The resparsifying streaming algorithm.

---

**Theorem 6.1** *Let  $\tilde{\mathbf{A}}$  be the matrix returned by  $\text{STREAMING-SAMPLE}(\mathbf{A}, \epsilon)$ . Then, with high probability,*

$$(1 - \epsilon)\mathbf{A}^\top \mathbf{A} \preceq \tilde{\mathbf{A}}^\top \tilde{\mathbf{A}} \preceq (1 + \epsilon)\mathbf{A}^\top \mathbf{A}.$$

**Sketch of proof:** To apply Theorem 3.1, it is enough to observe that, unless the algorithm already failed, the value  $l_j$  computed in step 3(b)i of  $\text{STREAMING-SAMPLE}$  is at most  $1 + \epsilon$  times smaller than the leverage score of  $\tilde{\mathbf{a}}_j$  wrt.  $\mathbf{A}$ . ■

## 7 Acknowledgments

The author would like to thank Richard Peng for suggesting the usage of sampling through unbiased coin flips, which significantly simplified the analysis. The author would also like to thank Michael Cohen and Yiannis Koutis for helpful comments and feedback.

## References

- [CLM<sup>+</sup>15] Michael B Cohen, Yin Tat Lee, Cameron Musco, Christopher Musco, Richard Peng, and Aaron Sidford. Uniform sampling for matrix approximation. In *Proceedings of the 6th Conference on Innovations in Theoretical Computer Science (ITCS)*, pages 181–190, 2015. 1
- [CMP16] M. B. Cohen, C. Musco, and J. Pachocki. Online Row Sampling. *ArXiv e-prints*, April 2016. 1
- [KL13] Jonathan A Kelner and Alex Levin. Spectral sparsification in the semi-streaming setting. *Theory of Computing Systems*, 53(2):243–262, 2013. (document), 1
- [KS16] R. Kyng and S. Sachdeva. Approximate Gaussian Elimination for Laplacians: Fast, Sparse, and Simple. *ArXiv e-prints*, May 2016. 1
- [LMP13] Mu Li, Gary L Miller, and Richard Peng. Iterative row sampling. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 127–136, 2013. 1
- [PS14] Richard Peng and Daniel A. Spielman. An efficient parallel solver for sdd linear systems. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing, STOC ’14*, pages 333–342, New York, NY, USA, 2014. ACM. 1
- [SS08] Daniel A. Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing (STOC)*, pages 563–568, 2008. (document), 1
- [Tro11] Joel Tropp. Freedman’s inequality for matrix martingales. *Electronic Communications in Probability*, 16:262–270, 2011. 2
- [Tro12] Joel A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434, 2012. 2